# On Data Mining in Inverse Scattering Problems: Neural Networks Applied to GPR Data Analysis

Salvatore Caorsi and Mattia Stasolla

Department of Electronics

University of Pavia

27100 – Pavia, Italy

Email: {name.surname}@unipv.it

*Abstract*—**This paper presents a (semi-)automatic processing technique for GPR data analysis. Exploiting the generalization capabilities of artificial neural networks (ANN), it will be shown that it is possible to feed a Multi-Layer Perceptron (MLP) with a suitable set of input features in order to determine the permittivity of a ground layer. A detailed performance assessment have proven that the algorithm provides very promising results, reconstructing with high accuracy the dielectric properties of both planar and rough surfaces. Some critical issues have anyway emerged that limit the effectiveness of the method to lossless media.**

*Index Terms*—**Inverse scattering, artificial neural networks, GPR.**

## I. INTRODUCTION

The term non-destructive testing (NDT) denotes a wide group of analysis techniques used in science and industry to evaluate the properties of a material, component or system without causing damage [1]. NDT covers a broad and interdisciplinary range of research fields, from biomedical engineering to geology, providing an excellent balance between quality and cost-effectiveness. A typical device used in NDT is the Ground Penetrating Radar (GPR), designed primarily for the detection of objects and/or interfaces below the earth's surface [2]. As a radar instrument, the GPR consists in a transmitting antenna (TX) which emits e.m. pulses towards the ground and a receiving antenna (RX) recording in a so-called *trace* the travel time and amplitude of the backscattered wavelets (see Fig. 1). In standard GPR measurements [3], the antennas are pulled along the survey track while traces are triggered at a fixed interval by a measurement wheel which is connected to the back of the antenna. This results in a series of traces which are finally displayed by the measurement software as a function of position and time in a so-called *radargram*, shown in Fig. 2. The analysis of radargrams requires long-term expertise, and no well-established techniques for an automatic processing can be still found. The scope of this paper is therefore to provide a (semi-)automatic processing chain for one of the open issues in NDT: the inversion of GPR data for the reconstruction of the dielectric characteristics of a surface layer. Such information is of great interest to many research areas, for instance when soil moisture content must be evaluated [4] or for measurements of crop canopy properties [5].
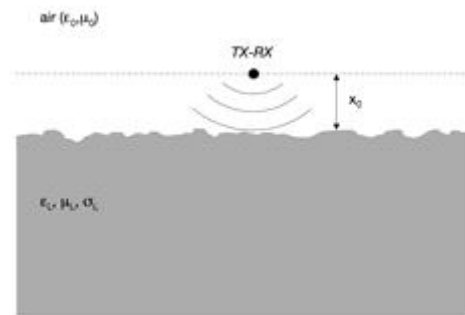


Figure 1. GPR system: a transmitting antenna (TX) emits e.m. pulses towards the ground, while a receiving antenna (RX) records the travel time and amplitude of the backscattered wavelets
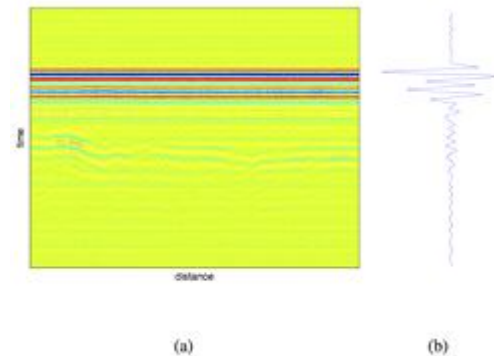


Figure 2. A radargram sample from on-field surveys (a). In (b) a single trace has been pointed out

To study in the most exhaustive way the problem, it must be first discussed an aspect which is often disregarded, i.e. how a rough surface can influence the received e.m. field. Commonly, solving approaches assume flat interfaces for simplification [6]-[7]; nevertheless, it has been demonstrated that realistic ground models can yield significant improvements [8]-[9]. In practice, in fact, the flat interface approximation could lead to appreciable errors, as shown in Fig. 3, where the received signal from different roughness profiles is depicted. It can be noticed that, with respect to the planar surface response in Fig.3a, the tail of signals can either exhibit ripples (Fig. 3b-c) or strong attenuation (Fig. 3d), due to multipath reflections. It is straightforward that such a different response could compromise the inversion process.
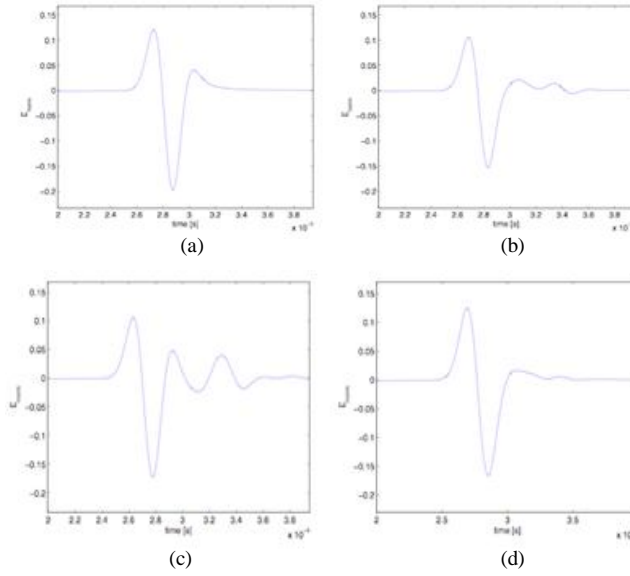
Figure 3. Effects of layer roughness on the received signals.

To solve the described problem, which would require a complex formalization of the physical phenomena involved, we here propose a method based on artificial neural networks (ANN). This kind of techniques have recently gained much interest in inverse scattering problems thanks to their capabilities of finding a regression model that relates the e.m. data to the desired outputs [10]. In the next section, the algorithm will be shown, introducing first some background information and then describing the actual processing scheme. Results and performance assessment will be presented in detail in Section III. Section IV concludes the discussion with final remarks.

## II. METHODOLOGY

Multi-Layer Perceptrons (MLPs) are feed forward artificial neural networks consisting of fully connected neurons arranged into layers, typically an *input*, *hidden* and *output* layer (see Fig. 4). Provided only that a proper number of hidden units and sufficiently smooth activation functions are available, they are capable of approximating any functional relation arbitrarily well [11]. In particular, for a three-layer MLP with topology *L-M-N*, the output $y_n$ of each output node is given by:

$$y_n = \sum_{m=1}^{M} f\left( \sum_{l=1}^{L} x_m w_{ml} + w_{m0} \right) w_{mn} + w_{n0} \qquad (1)$$

where $x$ is the input vector, $f$ represents a sigmoid function and $w_{ml}$ and $w_{nm}$ denote the input-hidden and hidden-output weights, which are internal adjustable parameters of the network. The knowledge of the MLP is acquired through a *supervised learning* process, where the availability of a a set

$$T = \left\{ \left( \mathbf{d_s}, y_s^* \right) \right\}_{s=1}^{S} \quad , \quad with \quad \mathbf{d_s} \in R^{1 \times M}$$

of ground-truthed input/output samples is assumed [12]. Given the *training set T*, the free parameters of the network

are adapted through a recursive minimization of the error

$$e = \sum_{s=1}^{S} \sum_{n=1}^{N} \frac{1}{2} \left( y_{sn} - y_{sn}^* \right)^2 \qquad (2)$$

between the actual and the expected values of the output, according to the Back Propagation (BP) algorithm [13]. Thanks to their generalization capabilities, MLPs seem therefore a suitable choice for reconstructing the permittivity of a layer by inverting GPR data. To this end, in agreement with the discussion above, a proper network topology and set of input features must be chosen. As the only parameter that needs to be reconstructed is the dielectric constant of the surface under test, the output layer of the MLP would be made of a single node. As far as the input layer is concerned, instead, we must consider that, provided that the *air wave* in bistatic GPRs can be windowed or assuming a monostatic GPR, the first wavelet reaching the RX is the *ground wave*, which is exactly the signal containing the information about layer's permittivity. Hence, we propose to sample the
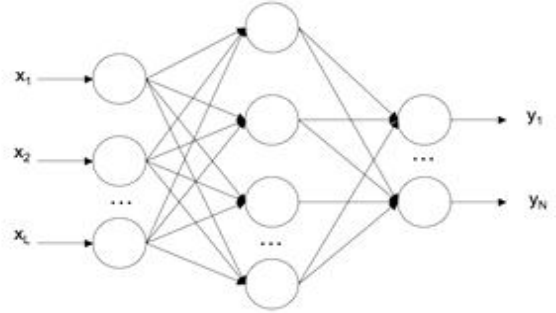


Figure 4. Architecture of a three-layer MLP, with *L* input nodes, *M* hidden nodes and *N* output nodes.

radargram trace at the times $t_k$ corresponding to its first $K$ maxima/minima $P_k$ and then feed the MLP with the vector:

$$\mathbf{d_s} = \left( P_{s1}, t_{s1}, P_{s2}, t_{s2}, \ldots, P_{sK}, t_{sK} \right)$$

The rationale behind this choice is that, on the one hand maxima and minima hold a higher informational content and, on the other hand, they can be quite simply detected through a peak detector. There are no well-established criteria to estimate which and how many input features, as well as the number of units within the hidden layer, should be employed to carry out the desired results. A rule of thumb states that, to ensure an adequate accuracy level, the appropriate number of training samples should not be less than the total number *W* of degrees of freedom (weights) of the network [14]:

$$S > W = M(L+1) + N(M+1) \qquad (3)$$

Therefore, the dimensioning of the network becomes a keypoint of the whole algorithm and should meet the need of a satisfying trade-off between network's complexity and its generalization capabilities. Noting that the transmitted signal has two local maxima and one minimum, and that any further peak implies the presence of multiple reflections, we opted to employ as input features the first *K=4* peaks, which would feed a *8-8-1 MLP*, thus with the same number of units for the input and hidden layer.
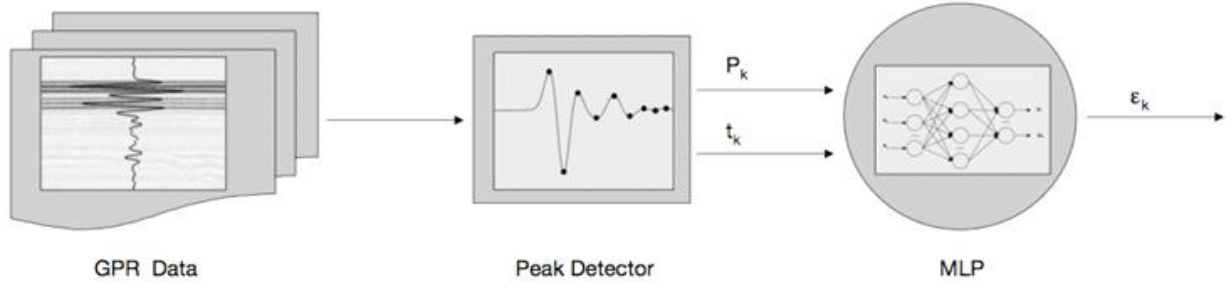
✦ACEEE

Figure 5. Algorithm workflow

The full processing scheme of the proposed algorithm for the NDT of surface layers is shown in Fig. 5. After the acquisition of GPR data, the radargram traces are processed by a peak detector, which extracts the amplitudes and corresponding times $\{P_k, t_k\}_{k=1}^4$ of the zero-derivative points of the signals. Then, the outputs of the peak detector are stored in a data vector and fed to the MLP, which finally returns the dielectric permittivity of the layer under test.

## III. RESULTS

The performance assessment of the algorithm described in the previous section has been carried out over a dataset of 200 different scenarios, simulated by means of GprMax, a software based on the FDTD numerical method [15].
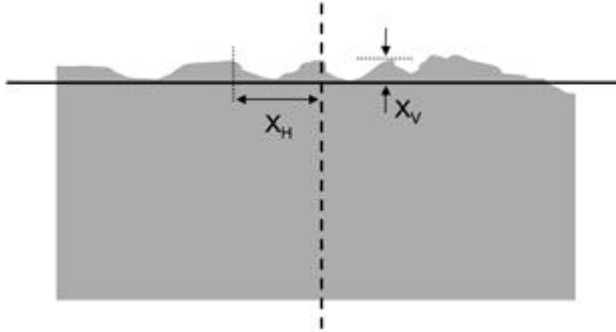


Figure 6.    Rough surface model: each discontinuity is generated imposing a vertical and horizontal shift from two reference planes.

To reduce computational costs, we have set a bi-dimensional simulation domain: thanks to the specific problem geometry, this approximation would not affect the validity of the test. According to the requirements of the algorithm, we have simulated a monostatic GPR, by modelling the TX as a current wire fed by a differentiated gaussian pulse of central frequency 2 GHz, and the RX as an ideal probe placed at the same position of the TX. As regards the geophysical properties of the layers, we have considered lossless media with permittivity $\varepsilon_L$ ranging from 2 to 20 and generated roughness profiles considering random vertical and horizontal deviations ($x_V$ and $x_H$) from suitable reference planes within the interval [0.5;1.5] cm (see Fig. 6). Such interval is consistent with the spatial resolution supplied by the system. The whole dataset has been first processed in order to extract the 8 input features needed by the MLP and then split into two separated groups employed as training and test set, respectively. It is interesting to notice that the four series of amplitudes $P_1$, $P_2$, $P_3$, $P_4$ –

depicted in Fig. 7 – do not exhibit the expected increasing monotone behavior (the higher the permittivity, the greater the backscattering), but fluctuate around a mean value. This implies that the same peak value of the electric field could be extracted in distinct signals corresponding to different permittivities (e.g. $P_3[\varepsilon_L=12]$ H' $P_3[\varepsilon_L=15]$), within the same signal at different positions (e.g. $P_3[\varepsilon_L=5]$ H' $P_4[\varepsilon_L=5]$), or both (e.g. $P_3[\varepsilon_L=5]$ H' $P_4[\varepsilon_L=7]$). Hence the importance of
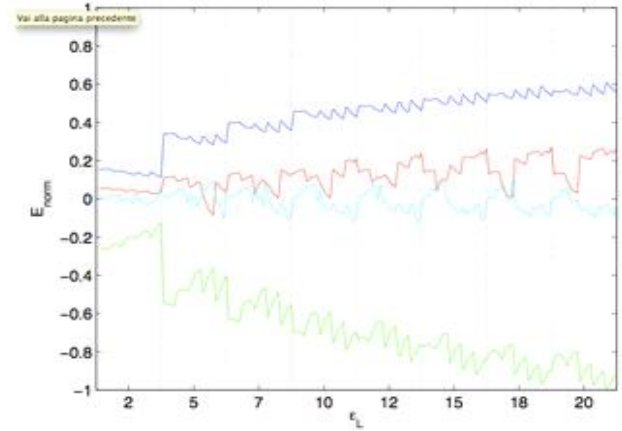


Figure 7. Electric field maxima/minima normalized values extracted within GPR data from different scenarios. P1: blue line; P2: green line; P3: red line; P4: cyan line.

TABLE I.
ABSOLUTE AND RELATIVE ERRORS FOR $S=100$.

| error | mean | max | min |
|---|---|---|---|
| relative | 2.50% | 9.32% | 0.00% |
| absolute | 0.28 | 1.69 | 0.00 |

feeding the MLP not only with amplitude values but also with the relative delays. According to Eq. 3 we expect that the network would show good generalization capabilities with more than 80 training samples. Table I shows the percentage relative error[1] between the network's output and the actual permittivity values in the case of $S=100$. It can be seen that the MLP performs very well, with a noticeable mean error of around 2.50% and a maximum error which remains below 10%. In absolute terms, it means that the dielectric constant of a layer can be reconstructed with an average accuracy of ±0.28, while the maximum spread should not exceed ±1.7. In addition, the graph in Fig. 8 illustrates that the mean error for each value of permittivity exhibits a satisfactorily quite uniform distribution of the error within the whole domain. Although the obtained results can be in general considered very promising, a training set of a hundred samples could, in

practice, become a limit for the effectiveness of the method. The main inconvenience in supervised approaches is, in fact, the necessity of finding a sufficient number of on-field examples in case no auxiliary data can be used for the learning phase. For this reason, we measured the network's performances by feeding it with an increasing number of I/O patterns, from 10 to 100, in order to evaluate the minimum dimension $S$ which could guarantee an adequate level of accuracy.
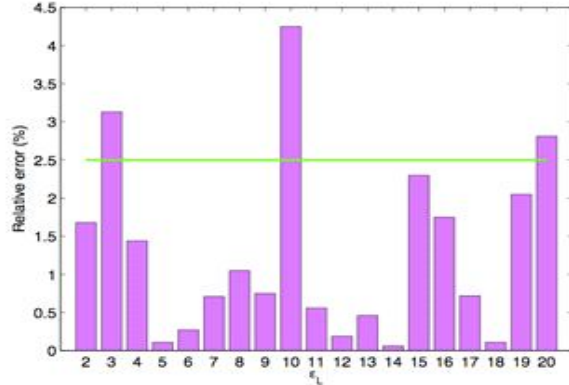


Figure 8. Mean relative error for each value of $\varepsilon_L$.

We actually found that the algorithm provides acceptable results also in case of very small training samples (see Tab II). Even in the worst case of $S=10$, it is indeed affected by a mean error of 14%, which is still satisfactory if we consider that it denotes an absolute mismatch with respect to the expected value of about 1.5. Maximum absolute errors seem, nevertheless, too high to justify the use of very small

training sets. A further aspect that should be mentioned regards the complementary trend shown by the error of the training and test set, which describe two opposite curves converging to a unique value (Fig. 9). This confirms that a proper fitting of the neural network yields good generalization performances (which exactly means that the expected error on unknown patterns should be very close the one obtained during the learning phase). To complete the performance assessment and test the method under additional realistic hypotheses, we evaluated the response of the MLP to e.m signals generated by planar surfaces and lossy media. In the former case, further simulations have been carried out employing the already mentioned horizontal reference plane as upper boundary of the layer. In the latter case, instead, we have generated both rough and planar media with a conductivity $s_L$ ranging from 0.01 to 10 S/m.
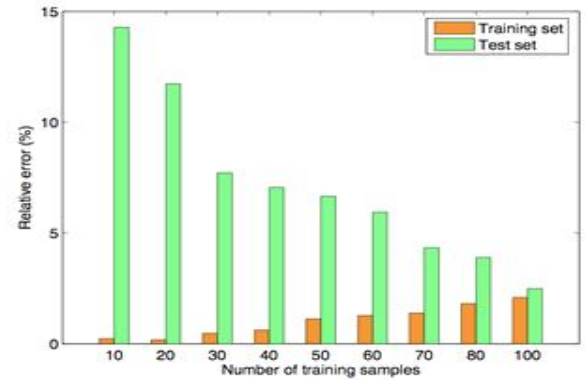


Figure 9. Mean relative error for increasing values of $S$.

TABLE II.
ABSOLUTE AND RELATIVE ERRORS FOR INCREASING VALUES OF $S$.

| | $S$ | 10 | | 20 | | 30 | | 40 | | 50 | | 60 | | 70 | | 80 | | 100 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | $\varepsilon$ | $\varepsilon_\%$ | $\varepsilon$ | $\varepsilon_\%$ | $\varepsilon$ | $\varepsilon_\%$ | $\varepsilon$ | $\varepsilon_\%$ | $\varepsilon$ | $\varepsilon_\%$ | $\varepsilon$ | $\varepsilon_\%$ | $\varepsilon$ | $\varepsilon_\%$ | $\varepsilon$ | $\varepsilon_\%$ | $\varepsilon$ | $\varepsilon_\%$ |
| **TRAIN** | max | 0.08 | 0.49 | 0.17 | 0.95 | 0.16 | 2.13 | 0.21 | 3.25 | 0.30 | 6.38 | 0.58 | 8.17 | 0.56 | 7.70 | 1.02 | 10.25 | 1.25 | 9.50 |
| | min | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 |
| | mean | 0.03 | 0.23 | 0.02 | 0.18 | 0.04 | 0.48 | 0.06 | 0.62 | 0.09 | 1.13 | 0.11 | 1.28 | 0.14 | 1.39 | 0.18 | 1.82 | 0.23 | 2.10 |
| **TEST** | max | 8.00 | 80.00 | 6.81 | 74.28 | 4.45 | 50.00 | 3.70 | 34.44 | 3.23 | 22.20 | 9.82 | 26.03 | 3.08 | 19.53 | 4,8 | 24.01 | 1.69 | 9.32 |
| | min | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 |
| | mean | 1.45 | 14.28 | 1.07 | 11.73 | 0.81 | 7.72 | 0.76 | 7.06 | 0.70 | 6.66 | 0.62 | 5.94 | 0.52 | 4.34 | 0.48 | 3.90 | 0.28 | 2.50 |

As illustrated in Fig. 10a, in case of planar surfaces network performances are comparable to those of the rough model (mean error of 2.69%), while in case of lossy layers the reconstruction is accomplished with a relative error of even 150%. Such behavior can be explained noticing that flat interface response is basically a particular case of rough response, so we can expect that the scattered signals would have amplitudes matching those employed within the training phase, and therefore properly invertible by the neural network. Lossy media, instead, suffer a twofold drawback: non-zero values of conductivity induce either a strong distortion of the backscattered signals, causing a considerable deformation of the wavelets (not recognizable by the MLP) or out-of-

domain amplitude values (not even manageable by the MLP). To improve performances it is possible to replace a few samples of the training set with input features extracted from lossy surface, and then re-train the MLP, in order to show the network how to recognize this kind of signals. A strong reduction of the errors can be indeed observed (see Fig. 10b), from 145% to 32%, even though, in absolute terms, the overall performances remain not completely satisfactory.
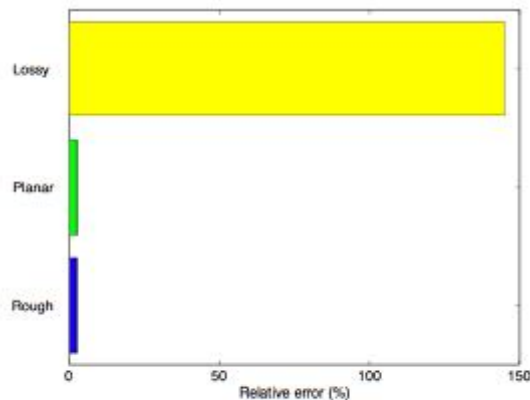
## IV. CONCLUSIONS

In this paper a (semi-)automatic ANN-based processing chain for GPR data analysis has been presented. In particular, we have shown that it is possible to feed a Multi-Layer Perceptron with a suitable set of input features extracted from the radargrams acquired during NDT in order to determine the permittivity of a surface layer. Supplied with a sufficient number of training samples, the network performs very well in presence of both planar and rough surfaces. Actually,

---

[1] the percentage relative error of the output $y$ with respect to the reference value $y^*$ has been computed according to the following formula:
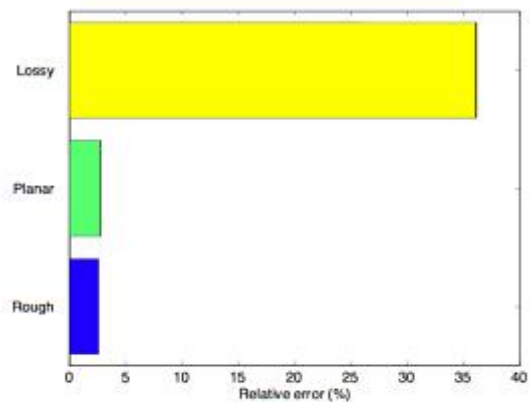
$$err(y) = \frac{|y^* - y|}{|y|} \times 100$$

satisfying generalization capabilities have been shown also when small training sets are provided. Such behavior suggests a possible over-parametrization of the neural network, whose design could be in case revised in order to reduce the overall complexity. A main drawback that has emerged from the analysis regards the limitations in handling signals scattered by lossy media: the distortion introduced by conductivity can induce wavelet's deformation and out-of-domain values, yielding to not completely satisfying performances.



(a)



(b)

Figure 10. Mean relative error for rough, planar and lossy media.

Future works will be therefore devoted to solve the above-mentioned open issues and will try to extend the proposed solution to multi-layered media. A particular concern will also regard the processing of noisy signals.

REFERENCES

[1] L. Cartz, Nondestructive Testing. ASM International, 1995.
[2] D. Daniels, D. J. Gunton, and H. F. Scott, "Introduction to subsurface radar," Proc. IEE, part F, vol. 135, pp. 278–320, Aug 1988.
[3] R. Ludwig, H. Gerhards, P. Klenk, U. Wollschl£ger, and J. Buchner, "Electromagnetic methods in applied geophysics," Institute of Environmental Physics Heidelberg University, 2009.
[4] C. Berthelot, D. Podborochynski, T. Saarenketo, B. Marjerison, and C. Prang, "Ground-penetrating radar evaluation of moisture and frost across typical saskatchewan road soils," Advances in Civil Engineering, 2010.
[5] G. Serbin and D. Or, "Ground-penetrating radar measurement of crop and surface water content dynamics," Remote Sensing of Environment, vol. 96, no. 1, pp. 119–134, 2005.
[6] C.-P. Kao, J. Li, Y. Wang, H. Xing, and C. R. Liu, "Measurement of Layer Thickness and Permittivity Using a New Multi-layer Model From GPR Data," Trans. on Geoscience and Remote Sensing, IEEE, vol. 45, no. 8, pp. 2463–2470, Aug 2007.
[7] P. Meincke, "Linear gpr inversion for lossy soil and a planar air-soil interface," Geoscience and Remote Sensing, IEEE Transactions on, vol. 39, no. 12, pp. 2713 –2721, dec 2001.
[8] L. Gürel and U. Oguz, "Simulations of ground-penetrating radars over lossy and heterogeneous grounds," IEEE Transactions on Geoscience and Remote Sensing, vol. 39, no. 6, pp. 1190–1197, jun 2001.
[9] N. Pinel, C. L. Bastard, L. Liu, C. Bourlier, and Y. Wang, "Rough thin pavement thickness estimation by GPR," in Proc. of IGARSS'09, Cape Town, South Africa, 11–17 Jul 2009.
[10] S. R. H. Hoole, "Artificial neural networks in the solution of inverse electromagnetic field problems," Trans. Magn., IEEE, vol. 29, no. 2, pp. 1931–1934, Mar 1993.
[11] K. Hornik, "Approximation capabilities of multilayer feedforward networks," Neural Networks, vol. 4, no. 2, pp. 251–257, 1991.
[12] S. Haykin, Neural Networks: A Comprehensive Foundation, 2nd ed. Prentice Hall, July 1998.
[13] S. A. Mutasem, B. O. Khairuddin, and A. N. Shahrul, "Back Propagation Algorithm: The Best Algorithm Among the Multi-layer Perceptron Algorithm," Int. Journal of Computer Science and Network Security, vol. 9, no. 4, pp. 378–383, Apr 2009.
[14] R. A. Dunne, "Multi-layer perceptron models for classification," Ph.D. Dissertation, Murdoch University, 2003.
[15] A. Giannopoulos, "Modelling Ground Penetrating Radar by GprMax," Construction Building Mater., vol. 19, no. 10, pp. 755–762, Dec 2005.

✣ACEEE